

# SEMANTICS-BASED CALORIE CALCULATOR

A Paper  
Submitted to the Graduate Faculty  
of the  
North Dakota State University  
of Agriculture and Applied Science

By  
Sravan Raghu Kumar Narra

In Partial Fulfillment of the Requirements  
for the Degree of  
MASTER OF SCIENCE

Major Department:  
Computer Science

March 2017

Fargo, North Dakota

North Dakota State University  
Graduate School

---

**Title**

SEMANTICS-BASED CALORIE CALCULATOR

---

**By**

Sravan Raghu Kumar Narra

---

The Supervisory Committee certifies that this *disquisition* complies with North Dakota State University's regulations and meets the accepted standards for the degree of

**MASTER OF SCIENCE**

SUPERVISORY COMMITTEE:

Dr. Juan Li

---

Chair

Dr. Jun Kong

---

Dr. Yarong Yang

---

Approved:

03/24/2017

---

Date

Brain M.Slator

---

Department Chair

## **ABSTRACT**

In recent years, people are considering healthy diet habits and many of them are trying to track and maintain their daily diet and consumption. To assist them, there are many applications available online and those applications are capable of recording calories for the ingredients consumed, but users must check individual calories and calculate total calories manually. In this paper, we propose a new technique to calculate calories for a given recipe in multiple formats. The new technique uses tokenization, hashing techniques and fuzzy matching for entity extraction and finally does the unit conversion to calculate calories. We compared the results of the proposed technique with the outcomes of the existing applications. These results proved that the new technique has the capacity to produce similar results compared to that of the existing applications and able to calculate calories for recipes in the different formats available on the internet.

## TABLE OF CONTENTS

ABSTRACT .....	iii
LIST OF FIGURES .....	vi
1. INTRODUCTION .....	1
2. RELATED WORK .....	3
3. METHODOLOGY .....	4
3.1. System Overview .....	4
3.2. Database .....	6
3.2.1. Design .....	6
3.2.2. Data Source.....	7
3.3. Calorie Calculation.....	7
3.3.1. Entity Extraction.....	7
3.3.1.1. Tokenization .....	7
3.3.1.2. Hashing Technique .....	10
3.3.1.3. Fuzzy Matching .....	13
3.3.2. Unit Conversion.....	14
3.4. Implementation.....	16
3.5. User Interface .....	16
3.5.1. Input Recipe.....	16
3.5.2. Calorie Calculation .....	17
3.5.3. Results .....	18
4. EVALUATION.....	20
4.1. Experimental Design .....	20

4.1.1. Dataset .....	20
4.1.2. Evaluation Metrics.....	20
4.1.2.1. Accuracy .....	21
4.1.2.2. Precision.....	21
4.1.2.3. Recall .....	21
4.1.2.4. F-Measure .....	22
4.2. Experimental Results.....	22
4.2.1. Comparisons .....	23
5. CONCLUSION AND FUTURE WORK .....	25
5.1. Conclusion.....	25
5.2. Future Work .....	25
6. REFERENCES .....	26

## LIST OF FIGURES

<u>Figure</u>	<u>Page</u>
1. Format Of A Recipe.....	4
2. System Flow Chart.....	5
3. Database ER Diagram.....	6
4. Tokenization .....	9
5. Multi Values For Single Hash Key.....	11
6. Entities With Common Ingredients .....	12
7. Entities Without Common Ingredients .....	13
8. Unit Conversion.....	15
9. Conversion Table Used [10].....	15
10. User Interface, Input Online Recipe .....	17
11. User Interface, Input Own Recipe .....	17
12. User Interface, Extraction And Unit Conversion.....	18
13. User Interface, Final Output .....	18
14. Evaluation, Basic Notations.....	21
15. Evaluation - Precision, Recall And F-Measure Scores.....	22
16. Comparison – Accuracy Based On Different Extraction Techniques. ....	23

## 1. INTRODUCTION

In recent years, people give a lot of thought to maintaining a healthy diet, considering healthy habits and staying fit. Most of the people are concentrating on knowing how many calories they consume on regular basis. As per [1] this news release article, American adults are choosing healthier foods, consuming healthier diets. The USDA research has shown that the diet quality has improved substantially from 2005 to 2010. Researchers found out that use of nutrition information, including the nutrition facts panel found on most of the food packages has increased in the recent years. Close to 50% of the adults are considering nutritional facts while making food choices.

Interestingly, per a google consumer survey made in Germany in October 2015 [2], 93% of people in Germany cook at least once a week, 63% of the people use the internet (search engines, websites) to get a recipe to cook. There are various reasons to choose an online recipe, as it gives us unlimited choices of recipes, can quickly find recipes in a specific category and they are for free. And the same survey stated that 67% of the people use either laptop or desktop to get recipes and 44% use their smartphone. These results indicate that most of the people are trying to get recipes and trying new recipes once a week. These results are also increasing rapidly.

In the blog [3], it was stated that the entity with the greatest influence on what Americans cook is not Costco or Trader Joe's, it's not the Food Network or The New York Times, it's Google. Every month, more than a billion searches made on google are for recipes. The recipes that the google search engine usually displays on its first page have a huge impact on what Americans cook. In an article by New York Times [4], it was stated that there were 10 million recipe searches made on a single day on google alone. This tells us that how people are interested in getting a recipe online rather than any other source.

Here comes the main problem, people want to stay healthy and want to cook as per the recipe found online. People want to know the calorie count for an online recipe. But, surprisingly, the websites where we can find online recipes do not show calorie count for the recipe, for example [5]. There are so many applications available online to know the number of calories in each ingredient. But, people don't prefer typing in each ingredient along with the serving sizes and know the calories in each ingredient used and calculate total calories for the entire recipe. There are some applications available to calculate calories for the entire recipe, but, there are certain limitations like the recipe should be only in a specific format.

During the scope of this paper, we propose a new technique to calculate calories for a given recipe in multiple formats. And at the end, we evaluated and compared results obtained by this new technique with some of the existing applications and it proved that the new technique can provide the same results and able to calculate calories for recipes in different formats.

## 2. RELATED WORK

The standard method for calculating the nutrition content like calories for a recipe is to chemically analyse the end outcome of the recipe once cooked [6]. But, it is not easy to perform chemical analysis of dishes since it involves higher costs (in terms of money and time). And this approach may need more set of assessments to perform. Considering the billions of recipes available online, this chemical analysis doesn't seem to be a practical solution to calculate calories.

As an alternative way, smart kitchen concept introduced by 'P. Chi, J. Chen, H. Chu and J. Lo' in [7] is a computer technology to improve home cooking by providing awareness on calories of food ingredients used while cooking. The smart kitchen uses sensors to track an ingredient and track the number of calories used to prepare a dish.

There are so many other approaches including the image recognition techniques to determine the calories from the pictures of the recipe after preparation [8]. These techniques first try to find out the main components based on pattern recognition and try to predict the calories based on the outcome. However, apart from the results that stated that ordinary people are showing interest in using it, the accuracy of these methods is low. And, the recipe needs to be cooked before to calculate recipes.

[9] Presents many improvised algorithms by considering the nutrition values lost through a cooking process that might vary the calories calculated before cooking and after cooking. These methods also need to have a recipe in a certain format to calculate calories. However, the experiments made by the authors proved that combining the calories for ingredients before cooking are acceptable results provided the ingredients are selected precisely.

In this paper, we work on extracting all the ingredients in a recipe and sum up the individual calories to provide the complete calories associated with a recipe.

### 3. METHODOLOGY

#### 3.1. System Overview

## Ingredients

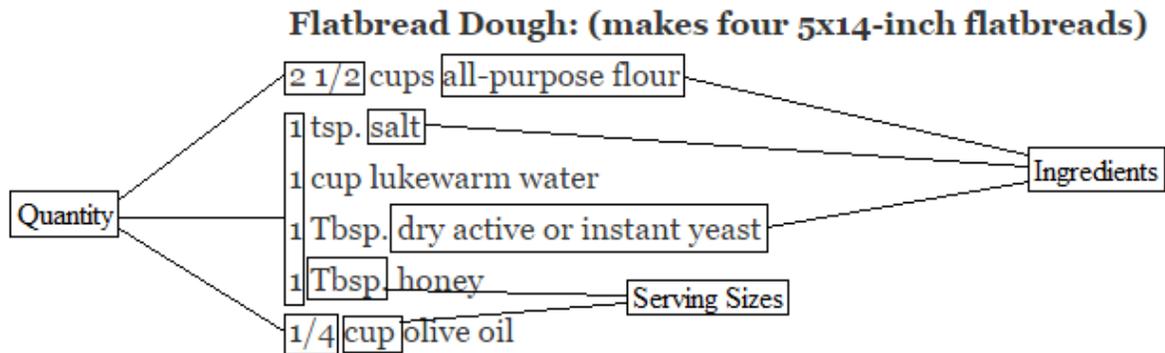


Figure 1. Format Of A Recipe

From the above figure, we can clearly determine how the recipe format looks like. An ideal recipe contains three parts for each ingredient used. They are the quantity of an ingredient used, serving size of an ingredient and the ingredient name. The recipes that are available online may or may not have quantity defined, may or may not have a serving size defined, but the ingredient name will be defined. The basic logic used to calculate the calories is presented in the flowchart below.

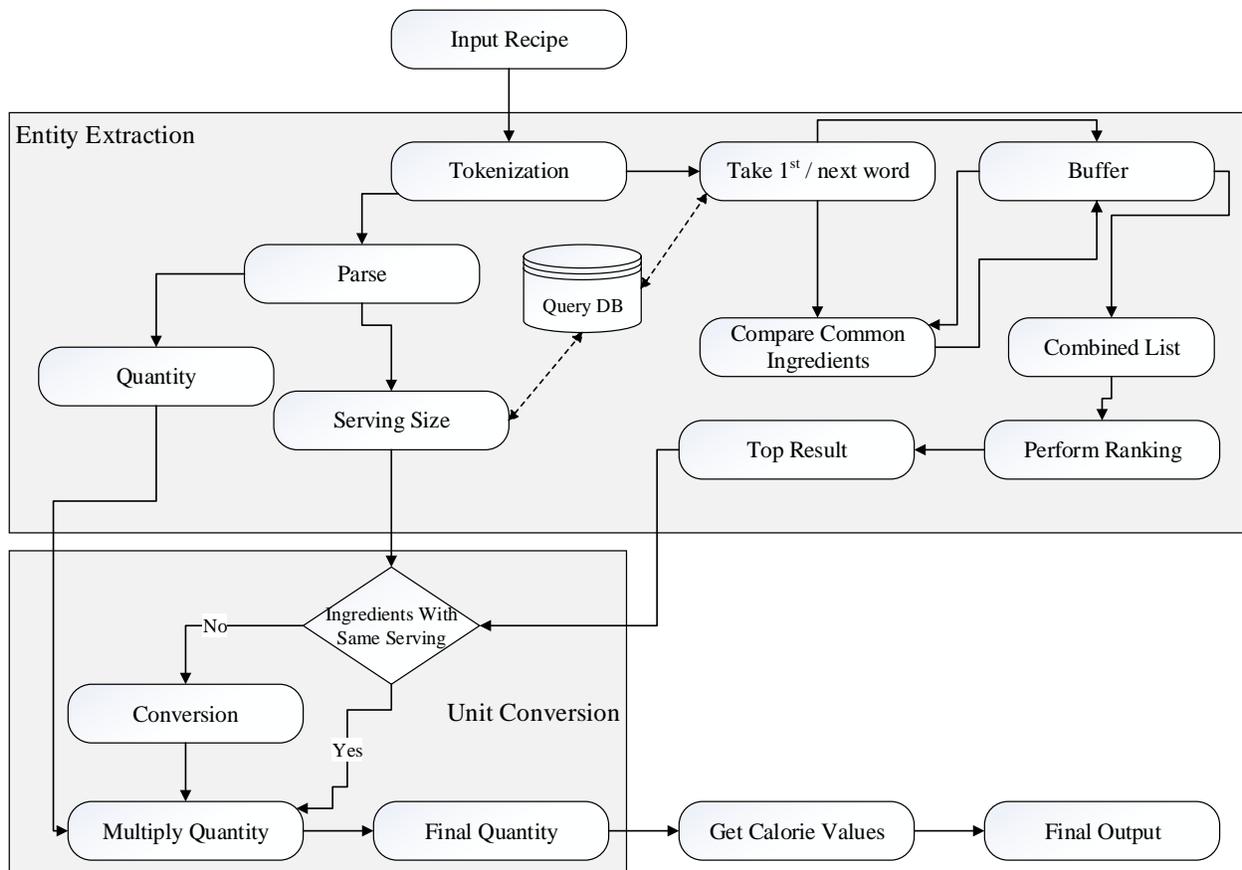


Figure 2. System Flow Chart

As presented in the above flowchart, the core idea contains two steps. Step 1 is termed as entity extraction, which is used to extract the given recipe into entities and classify those entities into three groups. Quantity group, serving group and ingredients group. In this step, we use tokenization to split the recipe into the smallest possible entities and to classify entities that belong to quantity and serving groups, then we use hash mapping technique to determine the ingredient, and finally, we use fuzzy matching to determine the ingredient name for those ingredients that are not figured out using hash mapping. By the end of this step, we can find out the quantities, servings and the ingredients used for the recipe. Step 2 is called as unit conversion, which is used to convert the quantity as per the serving sizes available for that ingredient. For this, we use a custom table that holds all the conversions for every serving size available. These are the 2 steps that are used

to calculate calories in every format available online. All these steps are explained later in this paper in detail.

### 3.2. Database

The main purpose of this database section is to get in details of how the database is designed considering all the techniques used to calculate calories and from where the initial data is loaded.

#### 3.2.1. Design

The figure below shows the ER Diagram for the whole database.

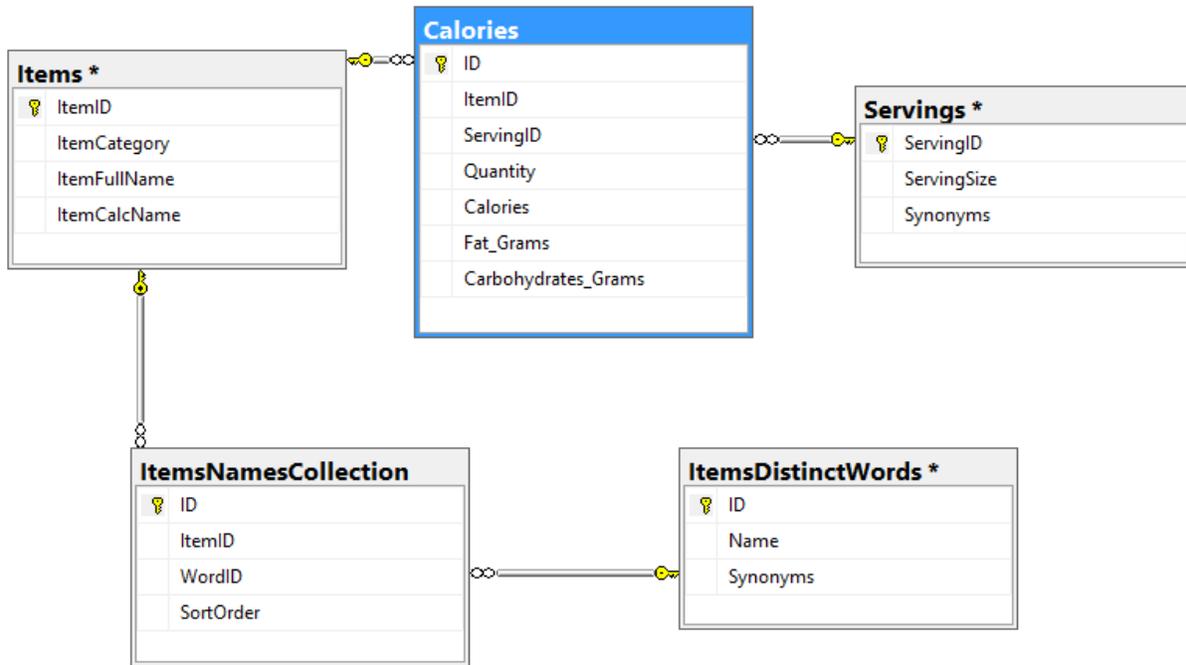


Figure 3. Database ER Diagram

We have used three base tables, 'Items' table to hold the ingredient names, 'Servings' table to hold the serving sizes and the core one, 'Calorie' table to hold the calorie information linking those other two tables. A couple of intermediate tables are used to make the extraction faster, one of them holds all the distinct words available in all the ingredients, and another one has a collection of ingredients that the word is part of. The basic idea to have these tables is to identify all the

ingredients for which an entity is available in. Apart from these five tables, we have used a function in the database to calculate the fuzzy match percentage between the ingredient from the recipe and the ingredients available in the database. For this reason, we have created a full-text search index in items table to retrieve the data faster.

### **3.2.2. Data Source**

There is no single place on the internet where the calorie data for all the ingredients is available. We have collected data from various sources so that we can cover most of the ingredients and calories. However, there is no approach that we can use to have calorie details for all ingredients. Here is a list of some of the sources we imported data from.

- What's Cooking America [31].
- Calorie Charts [32].
- Alberta Rose Calorie Charts [33].

## **3.3. Calorie Calculation**

The two steps used to calculate calories are 'Entity Extraction' and 'Unit Conversion'.

### **3.3.1. Entity Extraction**

Extraction is the concept of extracting the input recipe into smallest entities possible, and then classifying each entity as either quantity or serving or an ingredient. Entity extraction follows a three step process to extract and classify. The first stage is to use tokenization to split the input recipe into smallest possible entities.

#### ***3.3.1.1. Tokenization***

Tokenization is the process of breaking a stream of text up into words, phrases, symbols, or other meaningful elements called tokens [12]. Tokenization uses a tokenizer that breaks a text

or stream of text into tokens, usually by looking for a specified character. These are some of the steps that we go through tokenization.

The first step is used to split the input recipe into entities by removing all the spaces, extra characters like a newline (`\n`), tab (`\t`), etc. and punctuation marks. This will give us the list of all the entities used in the recipe.

Next step is to handle the abbreviations if there are any available in the input recipe. A dictionary is maintained with all the possible abbreviations, so that tokenization can look up for any word that is ending with a period in that dictionary and if found, it considers the dictionary result as an entity. For example, 'C.' is part of that dictionary that results to 'Cup'. When a recipe contains 'C.', tokenization knows that this is an abbreviation and consider 'Cup' as the resulting output. Most of the abbreviations are filtered if they are clearly defined, but there are some recipes where the abbreviations are not clearly defined. For example, 'lbs' is used instead of 'lbs.' [13], 'tbsp' is used instead of 'tbsp.' [13]. Such scenarios are also added in the dictionary and allowing tokenizer to look up for the words not ending with a period as well. The following figure shows some of the recipes that use this kind of notations.

Next step is to handle the numeric values or special expressions. This step determines if an entity belongs to quantity or not. If an entity is a numeric value, like '10' or '1/2', then we can determine it as quantity. But, there are some special expressions like '½' which cannot be classified as numeric by tokenizer [Ref]. We implemented similar dictionary concept to identify such type of scenarios so that tokenizer will result in the numeric values and classify that entity as quantity.

The following figure displays an example of tokenization process.

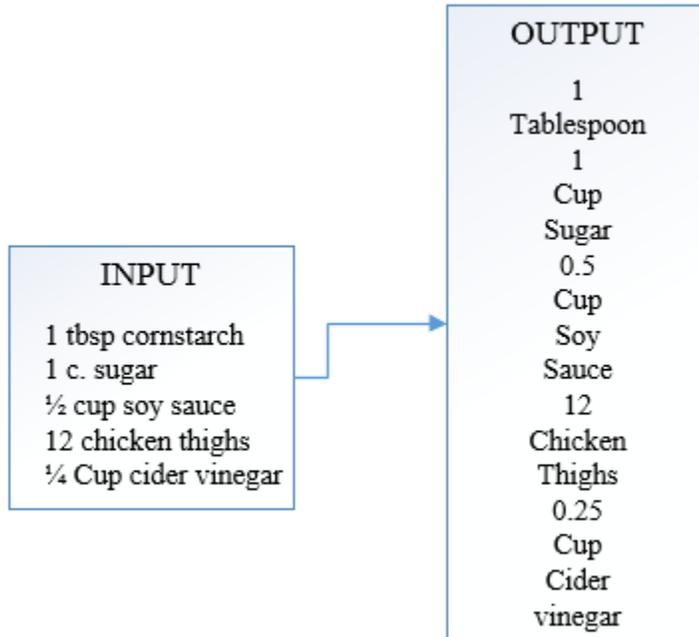


Figure 4. Tokenization

The resulting entities (tokens) of tokenization are used as an input for further processing. By the end of this stage, we will have all the entities used and the classification of quantities used in the recipe. The next stage is to classify all the servings available in the remaining entities.

To identify if an entity is a valid serving size or not, each entity is checked against the list of available servings in the database. Tokenization figures out all the abbreviations, synonym notations etc. We included some of the generic serving's sizes as well, like 'slice', 'medium', 'shank', 'inch' so that we can calculate calories accurately. By the end of this stage, we will have a classification of servings as well used in the recipe.

Every entity that was a result of tokenization process will be classified only into three groups. Quantity or serving size or an ingredient. If the entity does not belong to quantity and serving size, then that entity belongs to ingredients group. Next stage is to use hashing techniques

to determine whether that entity is a part of an ingredient like ‘sugar’ in ‘brown sugar’ or an ingredient itself like ‘salt’.

### ***3.3.1.2. Hashing Technique***

Hashing is a procedure that uniquely identifies a specific object from a group of similar objects. Hashing technique uses a hash table and a hash function. A hash table is a data structure that is used to store keys/value pairs [14]. For every Key, a value is stored so that retrieving the value is much faster. A hash table uses a hash function to compute an index into an array, from which the desired value can be found [34]. A hash function is used to distribute hash values uniformly. The concept of having a collection of Key, value pairs is also known as a map [15], and since hash function is used to determine the value of a key, it is known as a hash map. Multimap (sometimes also multihash) is a generalisation of a map in which more than one value may be associated with and returned for a given key [16]. To determine whether an entity is part of an ingredient like ‘sugar’ in ‘brown sugar’ or an ingredient itself like ‘salt’, we are using a similar technique to multihash.

The key to determining the exact ingredient is to consider all the entities irrespective of the classification groups, in the same order as resulted in tokenization. The following flow chart explains how to determine an ingredient from the entities.

Consider the first entity that does not belong to quantity group and servings group. Using multihash, we can identify all the ingredients that an entity is a part of (by passing entity as key, we can get a list of all the ingredients). The following figure shows an example of getting the result of multi-values for a single key (chocolate).

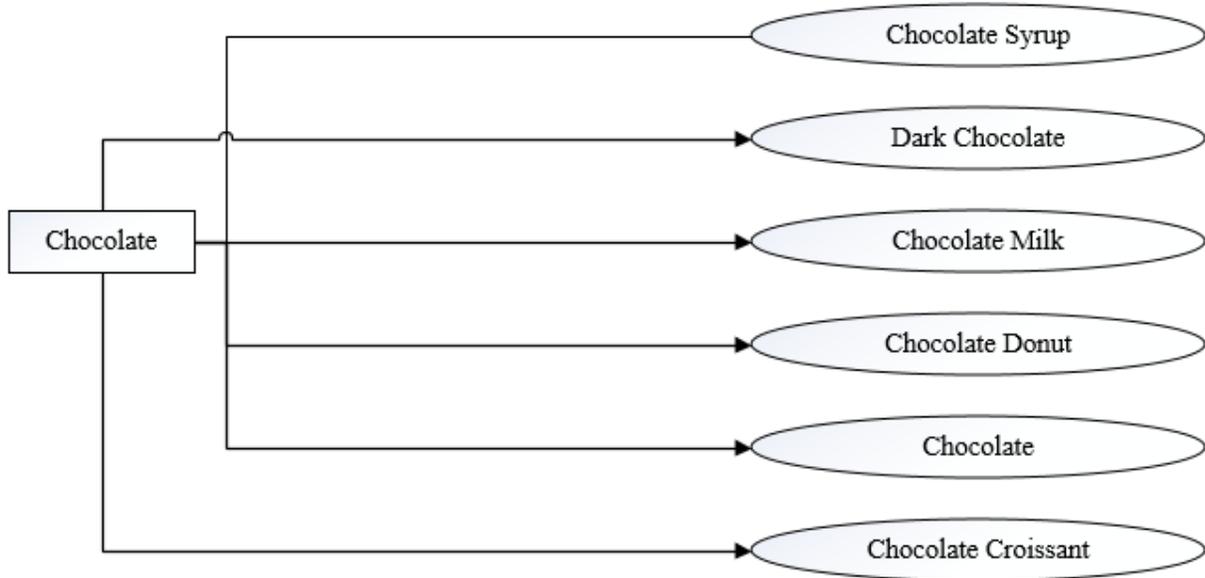


Figure 5. Multi Values For Single Hash Key

Then, if the next entity does not belong to quantity or serving size, we can check if there are any common ingredients between the ingredients resulted from the previous entity and the ingredients resulted from the current entity. If there are any common ingredients, then we proceed further with checking the common ingredients for next entity. The following figure shows an example for entities with common ingredients.

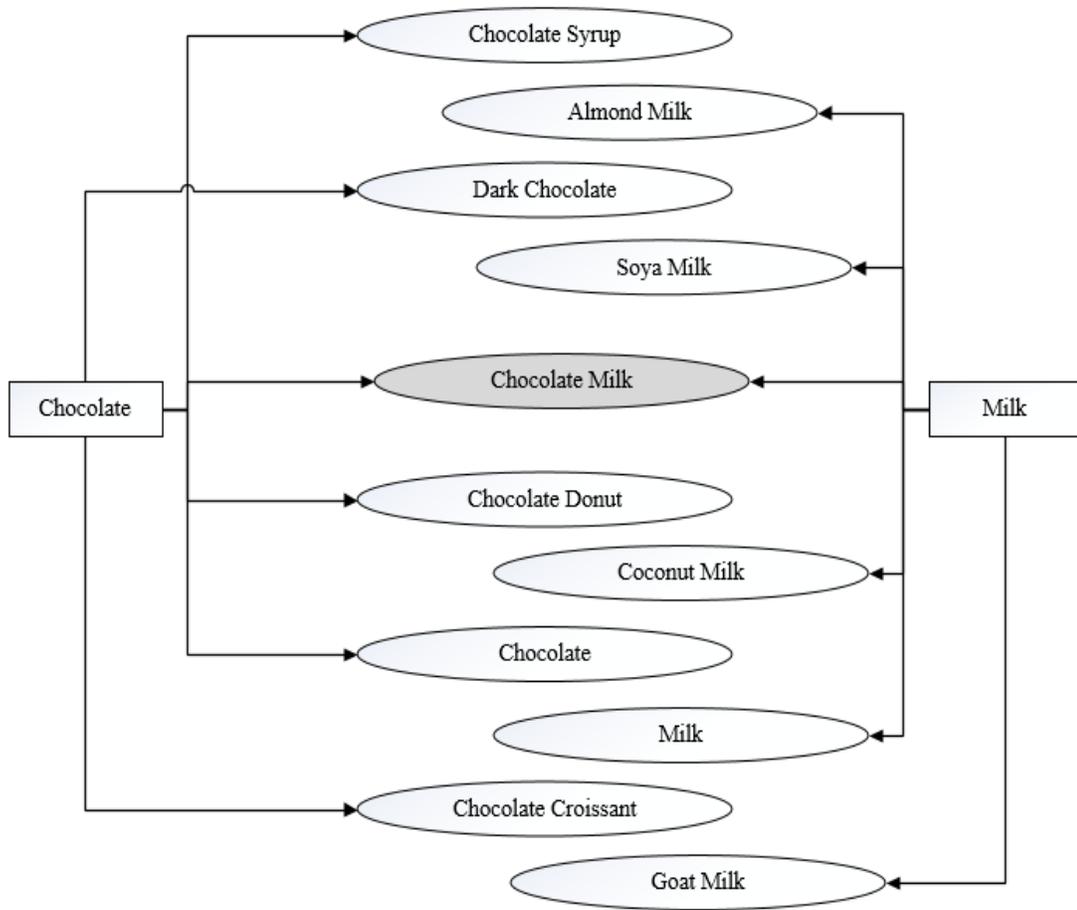


Figure 6. Entities With Common Ingredients

If there are no common entities between the previous entity ingredients and the current entity ingredients, it is an indication that the quantity and serving size are missing for the new ingredient. The following figure shows one such example.

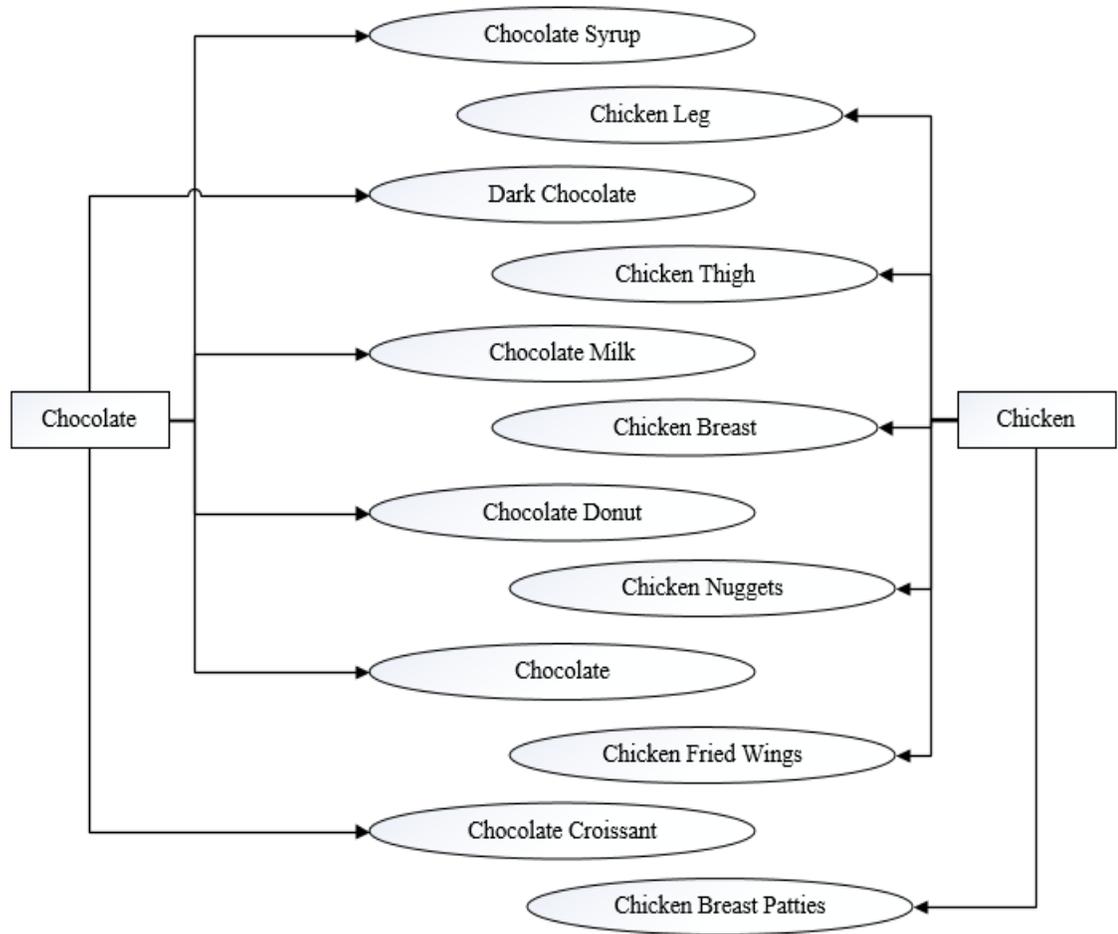


Figure 7. Entities Without Common Ingredients

If the next entity belongs to quantity or serving, then we can consider the matching entities as the ingredient framed from one or more entities.

There might be scenarios where the ingredient was typed incorrectly, as shown in the above figure. If this is the case, we combine all the concurrent unknown entities into individual ingredients and use fuzzy matching technique against all the ingredients available to identify these ingredients.

### 3.3.1.3. Fuzzy Matching

Fuzzy matching is a method that provides an improved ability to process word-based matching queries to find matching phrases or sentences from a database [17]. When there is no

100 percent match for a sentence or a phrase, fuzzy matching will try to find a match that's above a threshold matching percentage, which can be different based on different scenarios, which can be set by the application. We can implement fuzzy matching using many algorithms, but there is no single algorithm that was proved to be the best [18]. In this paper, authors researched seven different algorithms and concluded that none of the algorithms is the best for all values of the problem parameters.

A function in the database was created to determine the percentage match between any two strings. For all non-matching entities, we used this function to determine the percentage match between the ingredients from the recipe with all the ingredients available in the database. We are considering the ingredient that has the highest percentage match, only if the match percent is more than a threshold level. Currently, the threshold level was set to 80%. We considered different threshold levels like 90, 75, 65 etc. but 80% seems to yield better results.

If the exact ingredient name is figured out, then no need to follow through this phase. But, for some reason, if there are any difficulties determining the ingredient, we can use this phase to extract correct ingredient. By the end of this stage, we have extracted all the ingredients, quantities and serving sizes for all the ingredients used in the recipe.

### **3.3.2. Unit Conversion**

Unit conversion is the concept of converting the quantity for an ingredient from the serving size in the recipe matching with the serving size that we have calories for. For example, for an ingredient, the calories are saved for one cup, but the recipe is just using 3 tablespoons. For this, we need conversion from 1 cup to 1 tablespoon. To overcome such issues, we defined a custom method in our code that holds all the conversion from any serving size to all other serving sizes. Some of the examples are shown in the following figure.

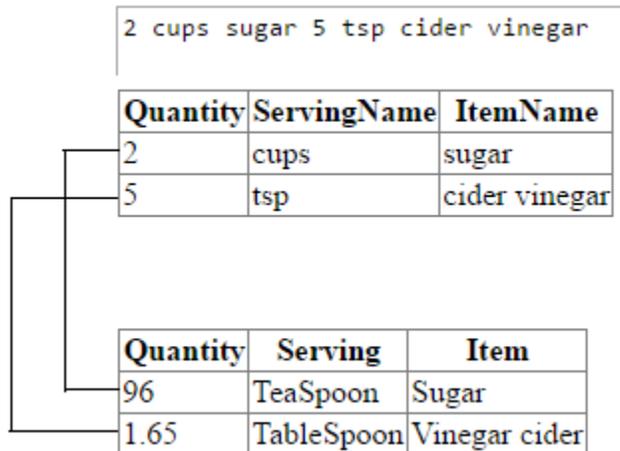


Figure 8. Unit Conversion

The following figure shows one example of the conversion tables used to convert to respective measures [10] [11].

### Measurements Conversion Chart

US Dry Volume Measurements	
MEASURE	EQUIVALENT
1/16 teaspoon	dash
1/8 teaspoon	a pinch
3 teaspoons	1 Tablespoon
1/8 cup	2 tablespoons (= 1 standard coffee scoop)
1/4 cup	4 Tablespoons
1/3 cup	5 Tablespoons plus 1 teaspoon
1/2 cup	8 Tablespoons
3/4 cup	12 Tablespoons
1 cup	16 Tablespoons
1 Pound	16 ounces

Figure 9. Conversion Table Used [10]

At the end of this phase, we will have the exact quantity, serving size and the exact ingredient to calculate calories.

There are some special scenarios where the ingredient used is not available in the database. If this is the case, we cannot calculate the calories correctly unless we know the calories for that ingredient. To overcome this issue, we added a functionality that a user can add ingredients and calories for that ingredient to the database. Once the ingredient is added, the user can recalculate the calories, which will show the exact calories in the recipe.

### **3.4. Implementation**

We implemented this as a web-based tool for which we need a user interface and we need a database that holds the calorie information and a program that is capable enough to connect to the database and retrieve data. We used SQL Server 2016 [19] to store the data and used .NET framework 4.5 [20] and C# language to create an interface. We can use any language, framework and any database until they are compatible with each other.

For tools, we used Visual Studio 2015 with C# to develop the application, and SQL Server Management Studio 2016 to query data.

### **3.5. User Interface**

The user interface contains three phases. Input Recipe, calorie calculation and Results.

#### **3.5.1. Input Recipe**

There are two ways the user can input the recipe. One is to use an online recipe and paste it in the application, other is to type in the recipe, for the recipes that were self-prepared.

If the user uses the recipe from online, the user can copy the recipe and paste it into the application. This type of input will go through entity extraction followed by unit conversion to calculate calories. The following figure shows an example recipe pasted from a website.

[Online Receipte](#)  
[Own Receipte](#)

```
4 eggs
2 cups sugar
1 teaspoon vanilla extract
2-1/4 cups all-purpose flour
2-1/4 teaspoons baking powder
1-1/4 cups 2% milk
10 tablespoons butter, cubed
```

Figure 10. User Interface, Input Online Recipe

If the user wishes to type in the recipe, the user can use ‘Own Recipe’ option and type in the recipe. This will not go through the entity extraction because we know the exact ingredients and serving sizes. However, this process needs to follow unit conversion to convert the quantity and get the correct serving sizes. The following figure shows an example recipe that was typed in by the user.

Quantity	Serving Size	Item Name	
<input type="text" value="4"/>	<input type="text" value="medium"/>	<input type="text" value="egg"/>	✗
<input type="text" value="2"/>	<input type="text" value="Cups"/>	<input type="text" value="sugar"/>	✗
<input type="text" value="1"/>	<input type="text" value="teaspoon"/>	<input type="text" value="van"/>	✗
<input type="text"/>	<input type="text"/>	vanilla extract	✗
<input type="text"/>	<input type="text"/>	almond milk, light vanilla	✗
<input type="text"/>	<input type="text"/>	shakes - vanilla	✗
<input type="button" value="Add New Row"/>			

[Online Receipte](#)  
[Own Receipte](#)

Figure 11. User Interface, Input Own Recipe

Once the recipe input is done, the user can click on ‘Calculate Calories’ to go through next stage and results in the output.

### 3.5.2. Calorie Calculation

This phase uses entity extraction if it was an online recipe and uses unit conversion for both input types to get the final result. We display the output of this phase to users to get a clear

idea of ingredients used in the recipe to the user. We display two tables, one is the outcome of the entity extraction phase and the other one is the outcome of unit conversion phase. The following figure shows both the outcomes for an example recipe.

Quantity	ServingName	ItemName
4	Each	eggs
2	cups	sugar
1	teaspoon	vanilla extract
2.25	cups	all purpose flour
2.25	teaspoons	baking powder
1.25	cups	2% milk
10	tablespoons	butter cubed

Quantity	Serving	Item
4	Each	Egg
96	TeaSpoon	Sugar
1	TeaSpoon	Vanilla Extract
2.25	Cup	All-Purpose Flour
2.25	TeaSpoon	Baking Powder
1.25	Cup	Milk, 2% fat
10	TableSpoon	Butter

Figure 12. User Interface, Extraction And Unit Conversion

### 3.5.3. Results

This is the final phase where the application shows the total calories for the entire recipe using the output of the previous phase. The following image shows the calories for an example recipe.

**Total Calories for the recipe are**

**4955**

Figure 13. User Interface, Final Output

The result might show like the calories may vary from AAA to BBB. This is due to the non-clear serving size. For example, if an onion is used and if the size of an onion is not mentioned we might show such results because a small onion will have different calories than a large onion.

We also have a small section in the user interface for the user to add the calories and ingredients for any that we do not have information.

## **4. EVALUATION**

### **4.1. Experimental Design**

In an experiment, we deliberately change one or more process variables (or factors) in order to observe the effect the changes have on one or more response variables [21]. In this paper, we are trying to evaluate the proposed technique to calculate calories for a specific set of data based on the evaluation metrics and compare the results with some of the existing websites that calculate calories.

#### **4.1.1. Dataset**

For this experiment, we considered 150 total recipes (10 different recipes from each website, from 15 different websites), some of them are [35] [36]. We combined all these 150 recipes into a single dataset and performed all the experiments on this dataset. All these recipes and websites that hold the recipes are randomly chosen. Performing an experiment on this dataset would be a perfect evaluation technique to determine the feasibility and accuracy of our algorithm, and compare the results with some of the existing applications.

#### **4.1.2. Evaluation Metrics**

There are so many metrics available to evaluate a technique [22] like Accuracy, Geometric Mean etc. For our experiment, we used four standard evaluation metrics. They are Accuracy, Precision, Recall and F-Measure to evaluate the performance. For all the evaluations, we used the following table as basic notations.

Total Population	Predicted Condition Positive	Predicted Condition Negative
Condition Positive	True Positive ( <i>tp</i> )	False Negative ( <i>fn</i> )
Condition Negative	False Positive ( <i>fp</i> )	True Negative ( <i>tn</i> )

Figure 14. Evaluation, Basic Notations

#### 4.1.2.1. Accuracy

The accuracy of a measurement system is the degree of closeness of measurements of a quantity to that quantity's true value [23]. Accuracy refers to the closeness of a measured value to a standard or a known value. Accuracy measures ratio of correct predictions to the total number of instances evaluated.

$$\text{Accuracy} = \frac{tp + tn}{(tp + fp + tn + fn)}$$

#### 4.1.2.2. Precision

Precision (also called positive predictive value) is the fraction of retrieved instances that are relevant [24]. Basically, precision answers the question ‘how many selected items are relevant.’

$$\text{Precision} = \frac{tp}{tp + fp}$$

#### 4.1.2.3. Recall

Recall (also known as sensitivity) is the fraction of relevant instances that are retrieved [24]. Basically, recall answers the question ‘how many relevant items are selected.’

$$\text{Recall} = \frac{tp}{tp + fn}$$

#### 4.1.2.4. F-Measure

F-measure is used to calculate F-score which is determined by calculating the harmonic mean of precision and recall.

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

## 4.2. Experimental Results

We have performed the proposed calorie calculation technique against the complete dataset with 150 recipes and recorded the precision, recall and F-Measure values for evaluation purpose.

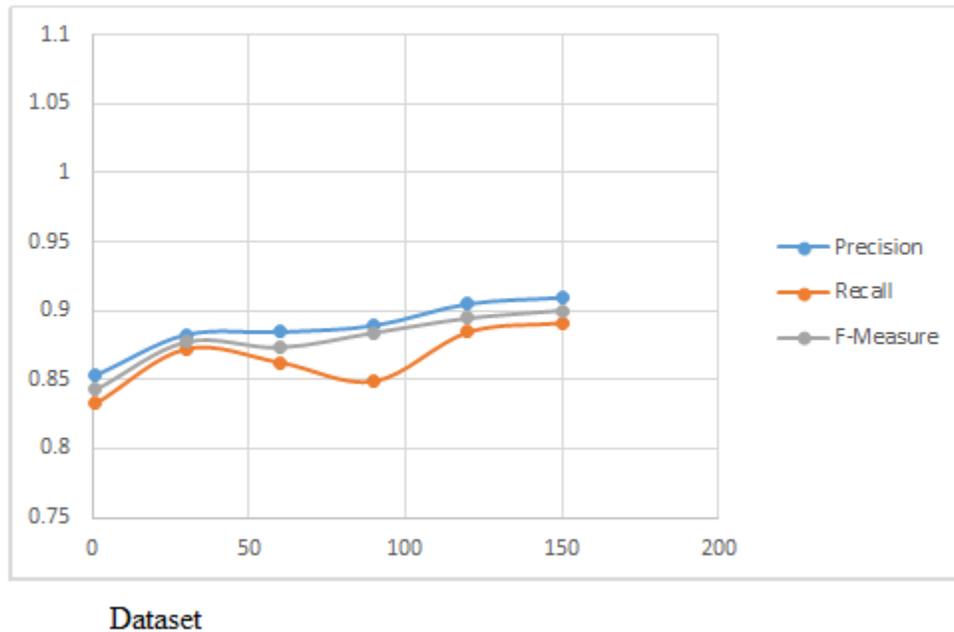


Figure 15. Evaluation - Precision, Recall And F-Measure Scores

The above graph shows precision, recall and F-Measure values that were achieved for the proposed technique. From the graph, we can see that the precision score is close to 91% and the recall is close to 89%. The better the precision and recall, the better the F-Measure scores. The F-Measure score is close to 90%. Now, let us compare this extraction results with other websites.

### 4.2.1. Comparisons

To compare the results, we are using two web applications [25], [26] and one mobile application [27] to calculate calories. There some other applications available to calculate the recipes, but they need step by step process to enter the recipe ingredients [28] [29] [30], like enter the first ingredient, enter the second ingredient etc. We cannot process the recorded dataset with these applications.

Even with the three applications considered, there are some restrictions like only one ingredient per line for all the applications. We cannot extract calories for the complete dataset with 150 recipes because all these applications are not designed to handle different formats of input recipe. So, we created a separate dataset that holds 50 recipes from various sources and compares the results for these 50 recipes.

We have recorded the accuracy for all the three websites for the new dataset with 50 recipes and compare these values with the proposed technique for evaluation purpose.

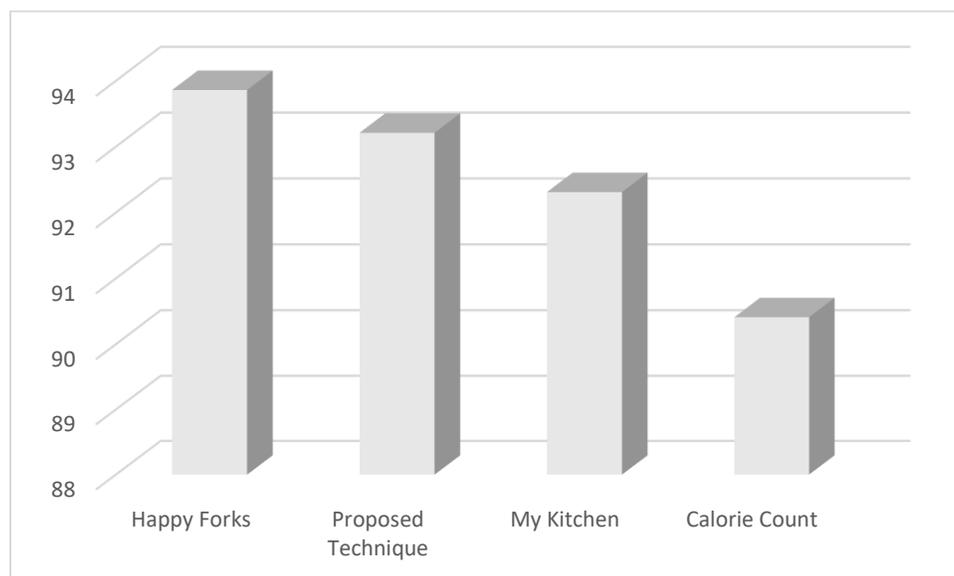


Figure 16. Comparison – Accuracy Based On Different Extraction Techniques.

From the above graph, we can evaluate the accuracy of calculating the calories for a given dataset with 50 recipes. It was seen that happyforks application resulted at 94%, mykitchen calculator at 92% and caloriecount resulted at 90%. The proposed technique resulted at 93%. All these differences in accuracy values should be the result of different extraction process that the applications are using, because all the calorie values will be same, and if all the applications extracted the same, then all the results will be the same. We are unable to comment on the application's extraction process since we do not know the extraction techniques used by the applications that we compared.

From the results of the experiment, it was seen that happyforks application recorded the highest accuracy at 94% and the proposed technique gave a tough competition for happyforks, recorded an accuracy at 93%. However, happyforks application and other applications are not capable of converting the recipes in a different format, which is the main purpose of the technique. Therefore, the proposed technique is capable enough to calculate calories in any format, with an accuracy of 93% and with the F-Measure score close to 90%.

## **5. CONCLUSION AND FUTURE WORK**

### **5.1. Conclusion**

In this paper, we addressed the problem of calculating calories for a recipe in different formats. For this, we have used tokenization, hashing techniques and fuzzy matching for entity extraction and used unit conversion for final conversions. And we have compared the results of the proposed technique with the existing applications to calculate calories for a given recipe. Experimental results on calorie calculation clearly show that the proposed technique is closely as accurate as the current applications and is capable of calculating calories for a recipe in different formats.

### **5.2. Future Work**

We can extend this to calculating other nutrition facts other than calories, like fat, carbohydrates etc. For this, we might need to explore more on the data collection part. And, we can use Natural Language Processing techniques in the extraction phase for better results but might take more time and big dataset for training. Hence, we can consider these improvements as the future work.

## 6. REFERENCES

- [1] American Adults are Choosing Healthier Foods, Consuming Healthier Diets (Jan 2014).  
Retrieved from  
<https://www.usda.gov/wps/portal/usda/usdahome?contentidonly=true&contentid=2014/01/0008.xml>
- [2] F. M. K. Gomez, and S. Neureuther (Nov 2015). The Role Of Internet For Cooking & Banking In Germany. Retrieved from  
[https://storage.googleapis.com/think-v2-emea/v2/a1721\\_2015%20Think%20With%20Google-%20Recipe%20Trends%20in%20Germany%20\(fmkg@%20-%20sneureuther@\)\\_v2.pptx.pdf](https://storage.googleapis.com/think-v2-emea/v2/a1721_2015%20Think%20With%20Google-%20Recipe%20Trends%20in%20Germany%20(fmkg@%20-%20sneureuther@)_v2.pptx.pdf)
- [3] Amanda and Merrill (May 2011). Googles New Recipe Search. Retrieved from  
<https://food52.com/blog/1838-update-google-s-new-recipe-search>
- [4] J. Moskin (May 2011). Can Recipe Search Engines Make You a Better Cook? Retrieved from  
<http://www.nytimes.com/2011/05/18/dining/can-recipe-search-engines-make-you-a-better-cook.html>
- [5] M. Collins (n.d.). Baked Teriyaki Chicken. Retrieved from  
<http://allrecipes.com/recipe/9023/baked-teriyaki-chicken/print/?recipeType=Recipe&servings=6>
- [6] Y Pico (2012). Chemical Analysis of Food: Techniques and Applications, Academic Book, University of Valencia.
- [7] P. Chi, J. Chen, H. Chu and J. Lo (June 2008). Enabling Calorie-Aware Cooking in a Smart Kitchen. Published in the 3rd international conference on Persuasive Technology.

- [8] K. Kitamura, C. de Silva, T. Yamasaki and K. Aizawa (July 2010). Image processing based approach to food balance analysis for personal food logging. Published in 2010 IEEE International Conference on Multimedia and Expo (ICME), pp. 625-630.
- [9] P. M. Powers and L. W. Hoover (March 1989). Calculating the nutrient composition of recipes with computers. *Journal of the American Dietetic Association* 89(2):224-32.
- [10] K. Maister (n.d.). Measurement and Conversion Charts. Retrieved from <http://startcooking.com/measurement-and-conversion-charts>
- [11] K. Isacks (n.d.). Useful Measurement Equivalents. Retrieved from <http://www.mynetdiary.com/estimating-portions-for-food-diary.html>
- [12] Tokenization (n.d.). Retrieved from [https://en.wikipedia.org/wiki/Tokenization\\_\(lexical\\_analysis\)](https://en.wikipedia.org/wiki/Tokenization_(lexical_analysis))
- [13] Ingredients for Chicken Teriyaki Recipe (Dec 2015). Retrieved from <http://natashaskitchen.com/2015/12/11/easy-teriyaki-chicken/>
- [14] P. Garg (n.d.). Basics of Hash Tables. Retrieved from <https://www.hackerearth.com/practice/data-structures/hash-tables/basics-of-hash-tables/tutorial/>
- [15] Associative Array (n.d.). Retrieved From [https://en.wikipedia.org/wiki/Associative\\_array](https://en.wikipedia.org/wiki/Associative_array)
- [16] Multimap (n.d.). Retrieved from <https://en.wikipedia.org/wiki/Multimap>
- [17] Fuzzy Matching (n.d.). Retrieved from <https://www.techopedia.com/definition/24183/fuzzy-matching>
- [18] P. Jokinen, J. Tarhio and E. Ukkonen (January 1988). A Comparison of Approximate String Matching Algorithms. *Software Practice and Experience*, Vol. 1(1), 1–4

- [19] SQL Server 2016 (n.d.). Retrieved from  
<https://www.microsoft.com/en-us/sql-server/sql-server-2016>
- [20] Microsoft .Net Framework 4.5 (n.d.). Retrieved from  
<https://www.microsoft.com/en-us/download/details.aspx?id=30653>
- [21] Experimental Design (n.d.). Retrieved from  
<http://www.itl.nist.gov/div898/handbook/pri/section1/pri111.htm>
- [22] M. Hossin and M.N. Sulaiman, A Review on Evaluation Metrics for Data Classification Evaluations. International Journal of Data Mining & Knowledge Management Process (IJDKP) Vol.5, No.2, March 2015.
- [23] Accuracy (n.d.). Retrieved from  
[https://en.wikipedia.org/wiki/Accuracy\\_and\\_precision](https://en.wikipedia.org/wiki/Accuracy_and_precision)
- [24] Precision and Recall (n.d.). Retrieved from  
[https://en.wikipedia.org/wiki/Precision\\_and\\_recall](https://en.wikipedia.org/wiki/Precision_and_recall)
- [25] Recipe Analyzer (n.d.). Retrieved from  
<https://happyforks.com/analyzer>
- [26] Recipe Analysis (n.d.). Retrieved from  
<https://www.verywell.com/recipe-nutrition-analyzer-4129594?ref=cc>
- [27] Recipe Converter (n.d.). Retrieved from  
<http://mykitchencalculator.com/recipeconverter.html>
- [28] Recipe Analyzer (n.d.). Retrieved from  
[https://www.eatracker.ca/recipe\\_analyzer.aspx](https://www.eatracker.ca/recipe_analyzer.aspx)
- [29] Recipe Calculator (n.d.). Retrieved from  
<https://recipes.sparkpeople.com/recipe-calculator.asp>

[30] Add Ingredient To Recipe (n.d.) Retrieved from  
[http://www.myfitnesspal.com/recipe/add\\_ingredient](http://www.myfitnesspal.com/recipe/add_ingredient)

[31] L. Stradley (n.d.). Food Nutrition Chart. Retrieved from  
<https://whatscookingamerica.net/NutritionalChart.htm>

[32] Calorie Charts (n.d.). Retrieved from  
<http://www.calorie-charts.net/>

[33] Calorie Charts (n.d.). Retrieved from  
[http://recipes.albertarose.org/calorie\\_charts/index.htm](http://recipes.albertarose.org/calorie_charts/index.htm)

[34] Hash Table (n.d.). Retrieved from  
[https://en.wikipedia.org/wiki/Hash\\_table](https://en.wikipedia.org/wiki/Hash_table)

[35] Cake Recipes (n.d.). Retrieved from  
<http://www.bettycrocker.com/recipes/dishes/cake-recipes>

[36] Food Recipes (n.d.). Retrieved from  
<http://www.food.com/>